

HIGH-INDEX OPTIMIZATION-BASED SHRINKING DIMER METHOD FOR FINDING HIGH-INDEX SADDLE POINTS*

JIANYUAN YIN[†], LEI ZHANG[‡], AND PINGWEN ZHANG[§]

Abstract. We present a high-index optimization-based shrinking dimer (HiOSD) method to compute index- k saddle points as a generalization of the optimization-based shrinking dimer method for index-1 saddle points [L. Zhang, Q. Du, and Z. Zheng, *SIAM J. Sci. Comput.*, 38 (2016), pp. A528–A544]. We first formulate a minimax problem for an index- k saddle point that is a local maximum on a k -dimensional manifold and a local minimum on its orthogonal complement. The k -dimensional maximal subspace is spanned by the k eigenvectors corresponding to the smallest k eigenvalues of the Hessian, which can be constructed by the simultaneous Rayleigh-quotient minimization technique or the locally optimal block preconditioned conjugate gradient method. Under the minimax framework, we implement the Barzilai–Borwein gradient method to speed up the convergence. We demonstrate the efficiency of the HiOSD method for computing high-index saddle points by applying finite-dimensional examples and semilinear elliptic problems.

Key words. rare event, saddle point, Morse index, dimer method, minimax

AMS subject classifications. 37M05, 49K35, 37N30, 34K28, 65P99

DOI. 10.1137/19M1253356

1. Introduction. The problems of finding saddle points on a potential energy surface (PES) have attracted much attention in different scientific communities during the past decades. A large spectrum of examples include finding critical nuclei and transition pathways in phase transformations [6, 35, 38, 41, 42], computing transition rates in chemical reactions [1, 25], and biology [37, 39, 44]. For instance, the transition state is characterized as an index-1 saddle point, that is, a critical point where the Hessian has one and only one negative eigenvalue. Because of the unstable nature of saddle points, methods for computing saddle points prove to be more challenging than those for a minimization problem.

A large number of numerical algorithms for computing index-1 saddle points have been proposed and applied to many practical problems. Generally speaking, there are two distinct approaches: path-finding methods and surface-walking methods. The former class requires initial and final states as two ends of the pathway for computing the minimum energy path (MEP) on the PES, and the index-1 saddle point corresponds to the configuration with the highest energy on the MEP. The representative methods include the nudged elastic band method [23, 24, 36] and the string method [8, 10, 11, 34]. The latter class starts from a single state on the PES and requires the first or the second derivatives of the potential energy to locate index-1 saddle points without a priori knowledge of the final state. Some samplers of the

*Submitted to the journal's Methods and Algorithms for Scientific Computing section March 29, 2019; accepted for publication (in revised form) September 24, 2019; published electronically November 12, 2019.

<https://doi.org/10.1137/19M1253356>

Funding: This work was supported by the National Natural Science Foundation of China through grants 11622102, 11421110001, 11861130351, 11421101, and 21790340. The first author is partially supported by the Elite Program of Computational and Applied Mathematics for PhD candidates in Peking University.

[†]School of Mathematical Sciences, Peking University, Beijing 100871, China (yinjy@pku.edu.cn).

[‡]Beijing International Center for Mathematical Research, Center for Quantitative Biology, Peking University, Beijing 100871, China (zhangl@math.pku.edu.cn).

[§]LMAM and School of Mathematical Sciences, Peking University, Beijing 100871, China (pzhang@pku.edu.cn).

surface-walking methods are the gentlest ascent dynamics [13, 15, 19], the dimer-type methods [17, 22, 23, 40, 43], the minimax method [27], the activation-relaxation technique [4, 29], the trajectory following algorithm [18], the step and slide method [31], the eigenvector-following method [5], the biased gradient squared descent method [9], etc. We refer to [12, 21, 45] as some excellent reviews.

While most existing algorithms are designed to find index-1 saddle points, the problems of computing high-index saddle points have received little attention despite its important applications. For instance, index-2 saddle points are particularly interesting in chemical systems for providing valuable information on the trajectories of chemical reactions [20]. The excited states in quantum mechanics that have higher energy than the ground states can also be characterized as high-index saddle point configurations because they could return to lower-energy excited states or ground states after excitation [14]. In terms of numerical algorithms, the path-finding methods can only find index-1 saddle points because the only unstable (maximum) direction is the direction along the MEP, and the other orthometric directions lead to the minimum. On the other hand, the surface-walking methods allow several approaches to find high-index saddle points. For instance, the minimax method [27] applied the local minimax theorem to compute index- k saddle points but required a priori knowledge of index- $(k - 1)$ saddle points. The biased gradient squared descent method [9] can also find high-index saddle points by transforming all critical points of the original potential energy into the minima of a gradient squared landscape, but the index of saddle points cannot be determined before calculation. The gentlest ascent dynamics (GAD) was developed by E and Zhou to search index-1 saddle points as a dynamical system of the gentlest ascent method [7], and such formulation can be extended to compute index-2 saddle points [13]. Quapp and Bofill developed a generalized GAD algorithm to locate high-index saddle points on the PES. It regards the saddle-searching problem as a dynamical system and requires the explicit calculation of the Hessian matrix [33].

Recently, Zhang, Du, and Zheng proposed the optimization-based shrinking dimer (OSD) method to efficiently compute index-1 saddle points [43], which can be regarded as a generalized formulation of the shrinking dimer dynamics [40]. The dimer rotation and translation steps in the dimer method [22] are transformed into corresponding optimization problems. Under the optimization framework, the OSD method shows a great advantage of applying optimization algorithms to accelerate the convergence. By adopting the similar spirit of the OSD, we present a high-index optimization-based shrinking dimer (HiOSD) method to search index- k saddle points. We first formulate the minimax problem for an index- k saddle point and then construct the maximal subspace by minimizing Rayleigh quotients simultaneously. Thus the HiOSD dynamical system is presented for finding an index- k saddle point, and the stability analysis is performed to show the linearly stable steady states of the HiOSD system are exactly index- k saddle points. To improve the numerical approximations of eigenvectors, we provide an alternative way by adopting the locally optimal block preconditioned conjugate gradient (LOBPCG) method to construct the maximal subspace. Furthermore, we apply the Barzilai–Borwein (BB) gradient method [3] to determine the step sizes in order to accelerate the convergence. Numerical examples, including finite-dimensional problems and semilinear elliptic problems, are presented to illustrate the efficiency of the HiOSD method for searching index- k saddle points.

2. HiOSD method. Given a twice Fréchet differentiable energy functional $E(\mathbf{x})$ defined on a real Hilbert space \mathcal{H} with an inner product $\langle \cdot, \cdot \rangle$, we let $\mathbf{F}(\mathbf{x}) = -\nabla E(\mathbf{x})$ denote its natural force and $\mathbb{G}(\mathbf{x}) = \nabla^2 E(\mathbf{x})$ denote its Hessian. By the Riesz

representation theorem, we regard $\mathbf{F}(\mathbf{x})$ as an element of \mathcal{H} . $\hat{\mathbf{x}} \in \mathcal{H}$ is called a *critical point* of $E(\mathbf{x})$ if $\|\mathbf{F}(\hat{\mathbf{x}})\| = 0$. A critical point of $E(\mathbf{x})$ that is not a local extremum is called a *saddle point* of $E(\mathbf{x})$. A critical point $\hat{\mathbf{x}}$ is called *nondegenerate* if $\mathbb{G}(\hat{\mathbf{x}})$ has a bounded inverse. According to the Morse theory [30], the *index* (Morse index) of a nondegenerate critical point $\hat{\mathbf{x}}$ is the maximal dimension of a subspace \mathcal{K} on which the operator $\mathbb{G}(\hat{\mathbf{x}})$ is negative definite. Our aim is to find an index- k saddle point, or for short, a *k-saddle*, on the PES. For simplicity, we assume the dimension of \mathcal{H} is d and write the inner product $\langle \mathbf{x}, \mathbf{y} \rangle$ as $\mathbf{x}^\top \mathbf{y}$.

2.1. Minimax optimization for a k -saddle. For a nondegenerate k -saddle $\hat{\mathbf{x}}$, the Hessian $\mathbb{G}(\hat{\mathbf{x}})$ has exactly k negative eigenvalues $\hat{\lambda}_1 \leq \dots \leq \hat{\lambda}_k$ with corresponding unit eigenvectors $\hat{\mathbf{v}}_1, \dots, \hat{\mathbf{v}}_k$ satisfying $\langle \hat{\mathbf{v}}_j, \hat{\mathbf{v}}_i \rangle = \delta_{ij}$, $1 \leq i, j \leq k$. By setting a k -dimensional subspace $\hat{\mathcal{V}} = \text{span}\{\hat{\mathbf{v}}_1, \dots, \hat{\mathbf{v}}_k\}$, $\hat{\mathbf{x}}$ is a local maximum on a k -dimensional linear manifold $\hat{\mathbf{x}} + \hat{\mathcal{V}}$ and a local minimum on $\hat{\mathbf{x}} + \hat{\mathcal{V}}^\perp$, where $\hat{\mathcal{V}}^\perp$ is the orthogonal complement space of $\hat{\mathcal{V}}$. By decomposing the Hilbert space \mathcal{H} as $\hat{\mathcal{V}} \oplus \hat{\mathcal{V}}^\perp$, the energy functional E can be defined with two inputs, respectively, in $\hat{\mathcal{V}}$ and $\hat{\mathcal{V}}^\perp$ naturally as $E(\mathbf{v}, \mathbf{w}) = E(\mathbf{v} + \mathbf{w})$, where $\mathbf{v} \in \hat{\mathcal{V}}$ and $\mathbf{w} \in \hat{\mathcal{V}}^\perp$. Therefore, the k -saddle $\hat{\mathbf{x}} = (\hat{\mathbf{x}}_{\hat{\mathcal{V}}}, \hat{\mathbf{x}}_{\hat{\mathcal{V}}^\perp})$ is a solution to the following minimax optimization problem:

$$(1) \quad \min_{\mathbf{w} \in \hat{\mathcal{V}}^\perp} \max_{\mathbf{v} \in \hat{\mathcal{V}}} E(\mathbf{v}, \mathbf{w}),$$

where $\hat{\mathbf{x}}_{\hat{\mathcal{V}}}$ is a local maximum of $E(\cdot, \hat{\mathbf{x}}_{\hat{\mathcal{V}}^\perp})$ on $\hat{\mathcal{V}}$, and $\hat{\mathbf{x}}_{\hat{\mathcal{V}}^\perp}$ is a local minimum of $E(\hat{\mathbf{x}}_{\hat{\mathcal{V}}}, \cdot)$ on $\hat{\mathcal{V}}^\perp$. Such a splitting structure exists around the nondegenerate k -saddle $\hat{\mathbf{x}}$ according to the Morse theory [27, 30]. It should be noticed that the subspace $\hat{\mathcal{V}}$ in (1) is unknown, which is an intrinsic difference from the classical minimax problems in optimization.

We first assume that \mathbf{x} is sufficiently close to $\hat{\mathbf{x}}$ so that $\mathbb{G}(\mathbf{x})$ has also exactly k negative eigenvalues. Let $\{(\lambda_i, \mathbf{v}_i)_{i=1}^k : \lambda_1 \leq \dots \leq \lambda_k < 0, \langle \mathbf{v}_j, \mathbf{v}_i \rangle = \delta_{ij}\}$ denote the eigenpairs of $\mathbb{G}(\mathbf{x})$, and we use a subspace $\mathcal{V} = \text{span}\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$ to approximate $\hat{\mathcal{V}}$ in (1). We define $\mathcal{P}_{\mathcal{V}}$ as the orthogonal projection operator on the finite-dimensional subspace $\mathcal{V} \subseteq \mathcal{H}$. To find a solution of the minimax problem (1), the dynamics of \mathbf{x} is supposed to satisfy that $\mathcal{P}_{\mathcal{V}}\dot{\mathbf{x}}$ is an ascent direction on the subspace \mathcal{V} and that $\dot{\mathbf{x}} - \mathcal{P}_{\mathcal{V}}\dot{\mathbf{x}}$ is a descent direction on the subspace \mathcal{V}^\perp . Therefore, we take $-\mathcal{P}_{\mathcal{V}}\mathbf{F}(\mathbf{x})$ as the direction of $\mathcal{P}_{\mathcal{V}}\dot{\mathbf{x}}$ and $\mathbf{F}(\mathbf{x}) - \mathcal{P}_{\mathcal{V}}\mathbf{F}(\mathbf{x})$ as the direction of $\dot{\mathbf{x}} - \mathcal{P}_{\mathcal{V}}\dot{\mathbf{x}}$, obtaining a gradient dynamics as

$$(2) \quad \dot{\mathbf{x}} = \beta_{\mathcal{V}}(-\mathcal{P}_{\mathcal{V}}\mathbf{F}(\mathbf{x})) + \beta_{\mathcal{V}^\perp}(\mathbf{F}(\mathbf{x}) - \mathcal{P}_{\mathcal{V}}\mathbf{F}(\mathbf{x})),$$

where $\beta_{\mathcal{V}}$ and $\beta_{\mathcal{V}^\perp}$ are positive relaxation constants. Alternatively, we may take a modified direction by choosing $\beta_{\mathcal{V}} = \beta_{\mathcal{V}^\perp} = \beta$ to get a modified gradient dynamics for a k -saddle as

$$(3) \quad \beta^{-1}\dot{\mathbf{x}} = \mathbf{F}(\mathbf{x}) - 2\mathcal{P}_{\mathcal{V}}\mathbf{F}(\mathbf{x}).$$

2.2. Construction of the subspace \mathcal{V} . The subspace \mathcal{V} in (3) is constructed by approximating the k eigenvectors simultaneously, and the approximate eigenvectors should also satisfy the orthonormal condition

$$(4) \quad \langle \mathbf{v}_i, \mathbf{v}_j \rangle = \delta_{ij}, \quad i, j = 1, 2, \dots, k.$$

The eigenvector \mathbf{v}_1 corresponding to the smallest eigenvalue λ_1 can be obtained by optimizing the Rayleigh quotient $\langle \mathbf{v}_1, \mathbb{G}(\mathbf{x})\mathbf{v}_1 \rangle / \langle \mathbf{v}_1, \mathbf{v}_1 \rangle$, or simply $\langle \mathbf{v}_1, \mathbb{G}(\mathbf{x})\mathbf{v}_1 \rangle$ with

a constraint $\langle \mathbf{v}_1, \mathbf{v}_1 \rangle = 1$. In a similar manner, computing the i th eigenvector \mathbf{v}_i can be transformed into a constrained optimization problem with the knowledge of $\mathbf{v}_1, \dots, \mathbf{v}_{i-1}$,

$$(5) \quad \min_{\mathbf{v}_i} \quad \langle \mathbf{v}_i, \mathbb{G}(\mathbf{x}) \mathbf{v}_i \rangle \quad \text{s.t.} \quad \langle \mathbf{v}_i, \mathbf{v}_j \rangle = \delta_{ij}, \quad j = 1, 2, \dots, i.$$

Then we minimize the k Rayleigh quotients (5) simultaneously [28] by using gradient-type dynamics of \mathbf{v}_i as follows:

$$(6) \quad \gamma^{-1} \dot{\mathbf{v}}_i = -\mathbb{G}(\mathbf{x}) \mathbf{v}_i + \langle \mathbf{v}_i, \mathbb{G}(\mathbf{x}) \mathbf{v}_i \rangle \mathbf{v}_i + \sum_{j=1}^{i-1} 2 \langle \mathbf{v}_j, \mathbb{G}(\mathbf{x}) \mathbf{v}_i \rangle \mathbf{v}_j, \quad i = 1, 2, \dots, k,$$

where $\gamma > 0$ is a relaxation parameter. To derive (6), the Lagrangian function of (5) is

$$(7) \quad \mathcal{L}_i(\mathbf{v}_i; \xi_1, \dots, \xi_{i-1}, \xi_i) = \langle \mathbf{v}_i, \mathbb{G}(\mathbf{x}) \mathbf{v}_i \rangle - \xi_i (\langle \mathbf{v}_i, \mathbf{v}_i \rangle - 1) - \sum_{j=1}^{i-1} \xi_j \langle \mathbf{v}_j, \mathbf{v}_i \rangle,$$

where ξ_1, \dots, ξ_i are Lagrangian multipliers, and the gradient of (7) is

$$(8) \quad \frac{\partial}{\partial \mathbf{v}_i} \mathcal{L}_i(\mathbf{v}_i; \xi_1, \dots, \xi_{i-1}, \xi_i) = 2\mathbb{G}(\mathbf{x}) \mathbf{v}_i - 2\xi_i \mathbf{v}_i - \sum_{j=1}^{i-1} \xi_j \mathbf{v}_j.$$

First, the dynamics for \mathbf{v}_1 is

$$(9) \quad \dot{\mathbf{v}}_1 = -\frac{\gamma}{2} \frac{\partial}{\partial \mathbf{v}_1} \mathcal{L}_1(\mathbf{v}_1; \xi_1) = -\gamma (\mathbb{G}(\mathbf{x}) \mathbf{v}_1 - \xi_1 \mathbf{v}_1).$$

Since $\langle \mathbf{v}_1, \dot{\mathbf{v}}_1 \rangle = 0$ should be satisfied from $\langle \mathbf{v}_1, \mathbf{v}_1 \rangle = 1$, we have $\xi_1 = \langle \mathbf{v}_1, \mathbb{G}(\mathbf{x}) \mathbf{v}_1 \rangle$, indicating the $i = 1$ case of (6). With $i \leq m - 1$ cases of (6), we set the dynamics of \mathbf{v}_m as

$$(10) \quad \dot{\mathbf{v}}_m = -\frac{\gamma}{2} \frac{\partial}{\partial \mathbf{v}_m} \mathcal{L}_m(\mathbf{v}_m; \xi_1, \dots, \xi_{m-1}, \xi_m) = -\gamma \left(\mathbb{G}(\mathbf{x}) \mathbf{v}_m - \xi_m \mathbf{v}_m - \frac{1}{2} \sum_{j=1}^{m-1} \xi_j \mathbf{v}_j \right).$$

$\xi_m = \langle \mathbf{v}_m, \mathbb{G}(\mathbf{x}) \mathbf{v}_m \rangle$ is derived from $\langle \mathbf{v}_m, \dot{\mathbf{v}}_m \rangle = 0$ similarly. From the condition $\langle \mathbf{v}_j, \mathbf{v}_m \rangle = 0$ for $j = 1, \dots, m - 1$, the dynamics (10) should satisfy $\langle \dot{\mathbf{v}}_j, \mathbf{v}_m \rangle + \langle \mathbf{v}_j, \dot{\mathbf{v}}_m \rangle = 0$, indicating $\xi_j = 4 \langle \mathbf{v}_j, \mathbb{G}(\mathbf{x}) \mathbf{v}_m \rangle$. As a consequence of introduction, we obtain the dynamics (6).

Under the orthonormal condition (4), the projection operator $\mathcal{P}_{\mathcal{V}}$ has a simple form of $\sum_{i=1}^k \mathbf{v}_i \mathbf{v}_i^\top$, and we obtain a dynamical system for a k -saddle:

$$(11) \quad \begin{cases} \beta^{-1} \dot{\mathbf{x}} = \left(\mathbb{I} - \sum_{i=1}^k 2\mathbf{v}_i \mathbf{v}_i^\top \right) \mathbf{F}(\mathbf{x}), \\ \gamma^{-1} \dot{\mathbf{v}}_i = - \left(\mathbb{I} - \mathbf{v}_i \mathbf{v}_i^\top - \sum_{j=1}^{i-1} 2\mathbf{v}_j \mathbf{v}_j^\top \right) \mathbb{G}(\mathbf{x}) \mathbf{v}_i, \quad i = 1, 2, \dots, k, \end{cases}$$

where \mathbb{I} is the identity operator and $\beta, \gamma > 0$ are relaxation parameters.

Because Hessians are often expensive to compute and store, we use first derivatives to approximate the Hessians in (11) by k shrinking dimers centered at \mathbf{x} , which is essentially central difference schemes for directional derivatives. The i th dimer ($i = 1, \dots, k$) has a direction of \mathbf{v}_i and length $2l$ and $\mathbb{G}(\mathbf{x})\mathbf{v}$ is approximated by

$$(12) \quad \mathbf{H}(\mathbf{x}, \mathbf{v}, l) = -\frac{\mathbf{F}(\mathbf{x} + l\mathbf{v}) - \mathbf{F}(\mathbf{x} - l\mathbf{v})}{2l}.$$

By adopting the idea of the shrinking dimer dynamics [40], we shrink the dimers by using a simple dynamics $\dot{l} = -l$, which corresponds to an exponential decay of the dimer length. In numerical algorithms, we set a minimum dimer length to avoid numerical errors due to round-off.

Now we obtain a complete HiOSD dynamical system for a k -saddle,

$$(13) \quad \begin{cases} \beta^{-1}\dot{\mathbf{x}} = \left(\mathbb{I} - \sum_{i=1}^k 2\mathbf{v}_i\mathbf{v}_i^\top \right) \mathbf{F}(\mathbf{x}), \\ \gamma^{-1}\dot{\mathbf{v}}_i = -\left(\mathbb{I} - \mathbf{v}_i\mathbf{v}_i^\top - \sum_{j=1}^{i-1} 2\mathbf{v}_j\mathbf{v}_j^\top \right) \mathbf{H}(\mathbf{x}, \mathbf{v}_i, l), \quad i = 1, 2, \dots, k, \\ \dot{l} = -l \end{cases}$$

with the orthonormal condition (4) satisfied initially.

Next we show that linearly stable steady states of (13) are exactly k -saddles.

THEOREM 1. *Assume that $E(\mathbf{x})$ is a \mathcal{C}^3 functional, $\mathbf{x}^* \in \mathcal{H}$, $\{\mathbf{v}_i^*\}_{i=1}^k \subset \mathcal{H}$ satisfies $\|\mathbf{v}_i^*\| = 1$, the Hessian $\mathbb{G}^* = \mathbb{G}(\mathbf{x}^*)$ is nondegenerate, whose eigenvalues are $\lambda_1^* < \dots < \lambda_k^* \leq \lambda_{k+1}^* \leq \dots \leq \lambda_d^*$, and $\beta, \gamma > 0$; then $(\mathbf{x}^*, \mathbf{v}_1^*, \dots, \mathbf{v}_k^*, l^*)$ is a linear stable steady state of (13) if and only if \mathbf{x}^* is a k -saddle of $E(\mathbf{x})$, \mathbf{v}_i^* is an eigenvector of \mathbb{G}^* corresponding to the eigenvalue λ_i^* , and l^* is 0.*

Proof. Taking the derivative of (12) with respect to l , we have

$$(14) \quad \frac{\partial}{\partial l} \mathbf{H}(\mathbf{x}, \mathbf{v}, l) = \frac{\mathbb{G}(\mathbf{x} + l\mathbf{v}) + \mathbb{G}(\mathbf{x} - l\mathbf{v})}{2l} \mathbf{v} + \frac{\mathbf{F}(\mathbf{x} + l\mathbf{v}) - \mathbf{F}(\mathbf{x} - l\mathbf{v})}{2l^2}.$$

As l decreases to zero, by the smoothness of $E(\mathbf{x})$ we have

$$(15) \quad \frac{\partial}{\partial l} \mathbf{H}(\mathbf{x}, \mathbf{v}, l) = \frac{\mathbb{G}(\mathbf{x})\mathbf{v} + o(l)}{l} + \frac{-2l\mathbb{G}(\mathbf{x})\mathbf{v} + o(l^2)}{2l^2} = o(1) \quad (\text{as } l \rightarrow 0).$$

We consider the Jacobian operator of the dynamics (13),

$$(16) \quad \mathbb{J} = \frac{\partial(\dot{\mathbf{x}}, \dot{\mathbf{v}}_1, \dot{\mathbf{v}}_2, \dots, \dot{\mathbf{v}}_k, \dot{l})}{\partial(\mathbf{x}, \mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k, l)} = \begin{pmatrix} \mathbb{J}_{\mathbf{x}} & \mathbb{J}_{\mathbf{x}1} & \mathbb{J}_{\mathbf{x}2} & \cdots & \mathbb{J}_{\mathbf{x}k} & \mathbb{O} \\ \mathbb{J}_{1\mathbf{x}} & \mathbb{J}_1 & \mathbb{O} & \cdots & \mathbb{O} & \mathbb{J}_{1l} \\ \mathbb{J}_{2\mathbf{x}} & \mathbb{J}_{21} & \mathbb{J}_2 & \cdots & \mathbb{O} & \mathbb{J}_{2l} \\ \vdots & \vdots & \vdots & & \vdots & \vdots \\ \mathbb{J}_{k\mathbf{x}} & \mathbb{J}_{k1} & \mathbb{J}_{k2} & \cdots & \mathbb{J}_k & \mathbb{J}_{kl} \\ \mathbb{O} & \mathbb{O} & \mathbb{O} & \cdots & \mathbb{O} & -1 \end{pmatrix},$$

whose blocks have following expressions:

$$(17) \quad \mathbb{J}_{\mathbf{x}} = \frac{\partial \dot{\mathbf{x}}}{\partial \mathbf{x}} = -\beta \left(\mathbb{I} - \sum_{i=1}^k 2\mathbf{v}_i \mathbf{v}_i^\top \right) \mathbb{G}(\mathbf{x}),$$

$$(18) \quad \mathbb{J}_{\mathbf{x}i} = \frac{\partial \dot{\mathbf{x}}}{\partial \mathbf{v}_i} = -2\beta (\mathbf{v}_i^\top \mathbf{F}(\mathbf{x}) \mathbb{I} + \mathbf{v}_i \mathbf{F}(\mathbf{x})^\top),$$

$$(19) \quad \mathbb{J}_i = \frac{\partial \dot{\mathbf{v}}_i}{\partial \mathbf{v}_i} = -\gamma \left(\mathbb{I} - \mathbf{v}_i \mathbf{v}_i^\top - \sum_{j=1}^{i-1} 2\mathbf{v}_j \mathbf{v}_j^\top \right) \frac{\mathbb{G}(\mathbf{x} + l\mathbf{v}) + \mathbb{G}(\mathbf{x} - l\mathbf{v})}{2} \\ + \gamma (\mathbf{v}_i^\top \mathbf{H}(\mathbf{x}, \mathbf{v}_i, l) \mathbb{I} + \mathbf{v}_i \mathbf{H}(\mathbf{x}, \mathbf{v}_i, l)^\top),$$

$$(20) \quad \mathbb{J}_{il} = \frac{\partial \dot{\mathbf{v}}_i}{\partial l} = -\gamma \left(\mathbb{I} - \mathbf{v}_i \mathbf{v}_i^\top - \sum_{j=1}^{i-1} 2\mathbf{v}_j \mathbf{v}_j^\top \right) \partial_l \mathbf{H}(\mathbf{x}, \mathbf{v}_i, l).$$

Furthermore, $\mathbb{J}_{\mathbf{x}i}$ is null if $\mathbf{F}(\mathbf{x}) = \mathbf{0}$, and \mathbb{J}_{il} is null if $l = 0$.

“ \Leftarrow ”: Supposing that \mathbf{x}^* is a k -saddle of $E(\mathbf{x})$, $(\lambda_i^*, \mathbf{v}_i^*)$ is an eigenpair of \mathbb{G}^* , and l^* is zero, we have $\mathbf{F}(\mathbf{x}^*) = \mathbf{0}$ and $\mathbf{H}(\mathbf{x}^*, \mathbf{v}_i^*, l^*) = \mathbb{G}^* \mathbf{v}_i^* = \lambda_i^* \mathbf{v}_i^*$, so $(\mathbf{x}^*, \mathbf{v}_1^*, \dots, \mathbf{v}_k^*, l^*)$ is a steady state of (13).

In order to show the linear stability, we consider the spectrum of $\mathbb{J}^* = \mathbb{J}(\mathbf{x}^*)$. From $\mathbf{F}(\mathbf{x}) = \mathbf{0}$ and $l = 0$, $\mathbb{J}_{\mathbf{x}i}$ and \mathbb{J}_{il} are null, so \mathbb{J}^* is block lower triangular. The first diagonal block,

$$(21) \quad \mathbb{J}_{\mathbf{x}}(\mathbf{x}^*, \mathbf{v}_1^*, \dots, \mathbf{v}_k^*, l^*) = \beta \left(-\mathbb{G}^* + \sum_{i=1}^k 2\lambda_i^* \mathbf{v}_i^* \mathbf{v}_i^{*\top} \right),$$

has eigenvalues of $\{\beta\lambda_1^*, \dots, \beta\lambda_k^*, -\beta\lambda_{k+1}^*, \dots, -\beta\lambda_d^*\}$, and the diagonal block,

$$(22) \quad \mathbb{J}_i(\mathbf{x}^*, \mathbf{v}_1^*, \dots, \mathbf{v}_k^*, l^*) = \gamma \left(-\mathbb{G}^* + \sum_{j=1}^i 2\lambda_j^* \mathbf{v}_j^* \mathbf{v}_j^{*\top} + \lambda_i^* \mathbb{I} \right),$$

has eigenvalues of $\{\gamma(\lambda_i^* + \lambda_1^*), \dots, \gamma(\lambda_i^* + \lambda_i^*), \gamma(\lambda_i^* - \lambda_{i+1}^*), \dots, \gamma(\lambda_i^* - \lambda_d^*)\}$. Because all eigenvalues of \mathbb{J}^* are negative, $(\mathbf{x}^*, \mathbf{v}_1^*, \dots, \mathbf{v}_k^*, l^*)$ is linear stable.

“ \Rightarrow ”: Supposing that $(\mathbf{x}^*, \mathbf{v}_1^*, \dots, \mathbf{v}_k^*, l^*)$ is a linear stable steady state, we have $l^* = 0$ straightforward, and therefore,

$$(23) \quad \left(\mathbb{I} - \mathbf{v}_i^* \mathbf{v}_i^{*\top} - \sum_{j=1}^{i-1} 2\mathbf{v}_j^* \mathbf{v}_j^{*\top} \right) \mathbb{G}^* \mathbf{v}_i^* = \mathbf{0}, \quad i = 1, 2, \dots, k.$$

Defining $\mu_i^* = \langle \mathbf{v}_i^*, \mathbb{G}^* \mathbf{v}_i^* \rangle$, we now show that for $i = 1, \dots, k$,

$$(24) \quad \mathbb{G}^* \mathbf{v}_i^* = \mu_i^* \mathbf{v}_i^*; \quad \langle \mathbf{v}_j^*, \mathbf{v}_i^* \rangle = \delta_{ij}, \quad j = 1, 2, \dots, i-1,$$

by introduction. Taking $i = 1$ in (23), we have $\mathbb{G}^* \mathbf{v}_1^* = \mu_1^* \mathbf{v}_1^*$ straightforward. Assuming that (24) holds for $i = 1, \dots, m-1$, by taking $i = m$ in (23), we have

$$(25) \quad \left(\mathbb{G}^* - \sum_{j=1}^{m-1} 2\mu_j^* \mathbf{v}_j^* \mathbf{v}_j^{*\top} \right) \mathbf{v}_m^* = \mu_m^* \mathbf{v}_m^*.$$

Since $\mathbf{v}_1^*, \dots, \mathbf{v}_{m-1}^*$ are eigenvectors of \mathbb{G}^* according to the inductive assumption, $\mathbb{G}^* - \sum_{j=1}^{m-1} 2\mu_j^* \mathbf{v}_j^* \mathbf{v}_j^{*\top}$ and \mathbb{G}^* share the same eigenvectors, so \mathbf{v}_m^* is an eigenvector of \mathbb{G}^* as well, and the eigenvalue is therefore μ_m^* by definition. From $\mathbb{G}^* \mathbf{v}_m^* = \mu_m^* \mathbf{v}_m^*$ and (25), we obtain $\sum_{j=1}^{m-1} \mu_j^* \langle \mathbf{v}_j^*, \mathbf{v}_m^* \rangle \mathbf{v}_j^* = 0$. Because $\{\mathbf{v}_j^*\}_{j=1}^{m-1}$ are orthogonal to each other, and the eigenvalue $\mu_j^* \neq 0$, we have $\langle \mathbf{v}_j^*, \mathbf{v}_m^* \rangle = 0$ for $j = 1, \dots, m-1$. As a consequence of introduction, we prove the property (24). Furthermore, we have $\mathbf{F}(\mathbf{x}^*) = \mathbf{0}$ from $(\mathbb{I} - \sum_{i=1}^k 2\mathbf{v}_i^* \mathbf{v}_i^{*\top})^2 = \mathbb{I}$, which indicates that \mathbf{x}^* is a critical point of $E(\mathbf{x})$.

Similarly, from $l^* = 0$ and $\mathbf{F}(\mathbf{x}^*) = \mathbf{0}$, \mathbb{J}^* is block lower triangular, so all eigenvalues of the diagonal blocks of \mathbb{J}^* are negative. We have proved that $\{(\mu_j^*, \mathbf{v}_j^*)\}_{j=1}^i$ are eigenpairs of \mathbb{G}^* , and then we let $\mu_{k+1}^* \leq \dots \leq \mu_d^*$ denote the other eigenvalues of \mathbb{G}^* . The eigenvalues of the first diagonal block,

$$(26) \quad \mathbb{J}_{\mathbf{x}^*, \mathbf{v}_1^*, \dots, \mathbf{v}_k^*, l^*} = \beta \left(-\mathbb{G}^* + \sum_{i=1}^k 2\mu_i^* \mathbf{v}_i^* \mathbf{v}_i^{*\top} \right),$$

are $\{\beta\mu_1^*, \dots, \beta\mu_k^*, -\beta\mu_{k+1}^*, \dots, -\beta\mu_d^*\}$, so $\{\mu_i^*\}_{i=k+1}^d$ are positive, and $\{\mu_i^*\}_{i=1}^k$ are negative, which means that \mathbf{x}^* is a k -saddle. The eigenvalues of the diagonal block,

$$(27) \quad \mathbb{J}_i(\mathbf{x}^*, \mathbf{v}_1^*, \dots, \mathbf{v}_k^*, l^*) = \gamma \left(-\mathbb{G}^* + \sum_{j=1}^i 2\mu_j^* \mathbf{v}_j^* \mathbf{v}_j^{*\top} + \mu_i^* \mathbb{I} \right),$$

are $\{\gamma(\mu_i^* + \mu_1^*), \dots, \gamma(\mu_i^* + \mu_i^*), \gamma(\mu_i^* - \mu_{i+1}^*), \dots, \gamma(\mu_i^* - \mu_d^*)\}$, so μ_i^* is smaller than μ_{i+1}^* for $i = 1, \dots, k$. Therefore, μ_i^* is exactly λ_i^* , which ends our proof. \square

We present Algorithm 1 as a discrete scheme for (13).

Algorithm 1. HiOSD method for a k -saddle.

Input: $k \in \mathbb{N}, l^{(0)} > 0, \mathbf{x}^{(0)} \in \mathcal{H}, \{\mathbf{v}_i^{(0)}\}_{i=1}^k \subset \mathcal{H}$ satisfying $\langle \mathbf{v}_j^{(0)}, \mathbf{v}_i^{(0)} \rangle = \delta_{ij}$.

- 1: Set $n = 0$, compute $\mathbf{f}^{(0)} = \mathbf{F}(\mathbf{x}^{(0)})$;
- 2: **repeat**
- 3: $\mathbf{x}^{(n+1)} = \mathbf{x}^{(n)} + \beta^{(n)} \mathbf{g}^{(n)}$;
- 4: **for** $i = 1 : k$ **do**
- 5: $\mathbf{v}_i^* = \mathbf{v}_i^{(n)} + \gamma_i^{(n)} \mathbf{d}_i^{(n)}$;
- 6: $\mathbf{v}_i^* = \mathbf{v}_i^* - \sum_{j=1}^{i-1} \langle \mathbf{v}_j^{(n+1)}, \mathbf{v}_i^* \rangle \mathbf{v}_j^{(n+1)}$;
- 7: $\mathbf{v}_i^{(n+1)} = \mathbf{v}_i^* / \|\mathbf{v}_i^*\|$;
- 8: **end for**
- 9: $l^{(n+1)} = \max \{l^{(n)} / (1 + \beta^{(n)}), \varepsilon\}$;
- 10: $\mathbf{f}^{(n+1)} = \mathbf{F}(\mathbf{x}^{(n+1)})$;
- 11: $n := n + 1$;
- 12: **until** $\|\mathbf{f}^{(n)}\| < \epsilon_F$;

Output: $\mathbf{x}^{(n)}, \mathbf{v}_1^{(n)}, \dots, \mathbf{v}_k^{(n)}$.

In Algorithm 1, $\mathbf{g}^{(n)}$ and $\mathbf{d}_i^{(n)}$ are defined as

$$(28) \quad \mathbf{g}^{(n)} = \mathbf{f}^{(n)} - 2 \sum_{i=1}^k \langle \mathbf{v}_i^{(n)}, \mathbf{f}^{(n)} \rangle \mathbf{v}_i^{(n)},$$

$$(29) \quad \mathbf{d}_i^{(n)} = -\mathbf{u}_i^{(n)} + \langle \mathbf{v}_i^{(n)}, \mathbf{u}_i^{(n)} \rangle \mathbf{v}_i^{(n)} + \sum_{j=1}^{i-1} 2 \langle \mathbf{v}_j^{(n)}, \mathbf{u}_i^{(n)} \rangle \mathbf{v}_j^{(n)},$$

$$(30) \quad \mathbf{u}_i^{(n)} = \mathbf{H}(\mathbf{x}^{(n+1)}, \mathbf{v}_i^{(n)}, l^{(n)}),$$

and $\beta^{(n)}, \gamma_i^{(n)}$ are the step sizes which will be determined in subsection 2.4. The Gram–Schmidt orthogonalization in step 6 and 7 is used to ensure that $\langle \mathbf{v}_j^{(n+1)}, \mathbf{v}_i^{(n+1)} \rangle = \delta_{ij}$ holds in discrete schemes. Dimer semilengths are lower-bounded by ε in step 9 for numerical stability, which is set as 10^{-6} in our numerical examples. Furthermore, to reduce the force evaluations, we can compute $\mathbf{f}^{(n+1)}$ by

$$(31) \quad \mathbf{f}^{(n+1)} = \frac{1}{k} \sum_{i=1}^k \frac{\mathbf{F}(\mathbf{x}^{(n+1)} + l^{(n)} \mathbf{v}_i^{(n)}) + \mathbf{F}(\mathbf{x}^{(n+1)} - l^{(n)} \mathbf{v}_i^{(n)})}{2}$$

in step 10, because the forces in the right side of (31) have been calculated in the previous iteration step.

2.3. Construction of subspace \mathcal{V} with LOBPCG. In discrete schemes, there are various efficient algorithms for finding eigenvectors of a large real symmetric matrix, but it is expensive and unnecessary to compute them accurately in each iteration step. To find better approximations of eigenvectors within admissible costs, we introduce the idea of the LOBPCG method [26].

The LOBPCG method is an iterative algorithm to find the smallest (or largest) k eigenvalues (or generalized eigenvalues) of a real symmetric matrix \mathbb{A} . Assuming $\{\mathbf{v}_i^{(n)}\}_{i=1}^k$ are approximations of k unit eigenvectors at the n th iteration step, we solve an optimization problem,

$$(32) \quad \min \sum_{i=1}^k \langle \mathbf{v}_i, \mathbb{A} \mathbf{v}_i \rangle \quad \text{s.t. } \mathbf{v}_i \in \mathcal{U}^{(n)}, \langle \mathbf{v}_j, \mathbf{v}_i \rangle = \delta_{ij},$$

exactly with Rayleigh–Ritz methods to obtain $\{\mathbf{v}_i^{(n+1)}\}_{i=1}^k$, where the subspace $\mathcal{U}^{(n)}$ is constructed according to $\{\mathbf{v}_i^{(n)}\}_{i=1}^k$. Given a symmetric positive definite preconditioner \mathbb{T} , the locally optimal block preconditioned steepest descent (LOBPSD) method chooses a subspace

$$(33) \quad \mathcal{U}_{\text{SD}}^{(n)} = \text{span} \left\{ \mathbf{v}_i^{(n)}, \mathbb{T} \left(\mathbb{A} \mathbf{v}_i^{(n)} - \langle \mathbf{v}_i^{(n)}, \mathbb{A} \mathbf{v}_i^{(n)} \rangle \mathbf{v}_i^{(n)} \right) : i = 1, \dots, k \right\},$$

and the LOBPCG method chooses a larger subspace

$$(34) \quad \mathcal{U}_{\text{CG}}^{(n)} = \text{span} \left\{ \mathbf{v}_i^{(n-1)}, \mathbf{v}_i^{(n)}, \mathbb{T} \left(\mathbb{A} \mathbf{v}_i^{(n)} - \langle \mathbf{v}_i^{(n)}, \mathbb{A} \mathbf{v}_i^{(n)} \rangle \mathbf{v}_i^{(n)} \right) : i = 1, \dots, k \right\}.$$

Since the matrix $\mathbb{A} = \mathbb{G}(\mathbf{x})$ is different at each iteration step and the terms with Hessians are approximated by dimers, we dilate on how the LOBPSD method is

applied in the HiOSD method. We first apply the LOBPSD method in each iteration step and use dimers (12) to approximate $\mathbb{G}(\mathbf{x})\mathbf{v}$ as well. We define $\mathbf{u}_i^{(n)}$ as (30), and

$$(35) \quad \mathbf{w}_i^{(n)} = \mathbb{T} \left(\mathbf{u}_i^{(n)} - \left\langle \mathbf{v}_i^{(n)}, \mathbf{u}_i^{(n)} \right\rangle \mathbf{v}_i^{(n)} \right), \quad i = 1, \dots, k,$$

and we compute the orthogonal basis of the subspace $\mathcal{U}_{\text{SD}}^{(n)}$ based on the normal orthogonal set $\{\mathbf{v}_i^{(n)}\}_{i=1}^k$. Specifically, we calculate

$$(36) \quad \tilde{\mathbf{w}}_i^{(n)} = \mathbf{w}_i^{(n)} - \sum_{j=1}^k \left\langle \mathbf{w}_i^{(n)}, \mathbf{v}_j^{(n)} \right\rangle \mathbf{v}_j^{(n)} - \sum_{\substack{j=1 \\ \|\tilde{\mathbf{w}}_j^{(n)}\| > \epsilon_w}}^{i-1} \frac{\left\langle \mathbf{w}_i^{(n)}, \tilde{\mathbf{w}}_j^{(n)} \right\rangle}{\|\tilde{\mathbf{w}}_j^{(n)}\|^2} \tilde{\mathbf{w}}_j^{(n)},$$

using the method of (modified) Gram–Schmidt recursively, and define a K -column matrix ($k \leq K \leq 2k$)

$$(37) \quad \mathbb{U}_{\text{SD}}^{(n)} = \left[\mathbf{v}_1^{(n)}, \dots, \mathbf{v}_k^{(n)}, \tilde{\mathbf{w}}_i^{(n)} / \|\tilde{\mathbf{w}}_i^{(n)}\| : \|\tilde{\mathbf{w}}_i^{(n)}\| > \epsilon_w, i = 1, \dots, k \right],$$

whose column vectors are orthogonal basis of the subspace $\mathcal{U}_{\text{SD}}^{(n)}$ approximately. In (37), we drop $\mathbf{w}_i^{(n)}$ whose norm is small enough for numerical stability. We define another K -column matrix

$$(38) \quad \mathbb{Y}_{\text{SD}}^{(n)} = \left[\mathbf{u}_1^{(n)}, \dots, \mathbf{u}_k^{(n)}, \mathbf{y}_i^{(n)} : \|\tilde{\mathbf{w}}_i^{(n)}\| > \epsilon_w, i = 1, \dots, k \right],$$

where

$$(39) \quad \mathbf{y}_i^{(n)} = \mathbf{H} \left(\mathbf{x}^{(n+1)}, \tilde{\mathbf{w}}_i^{(n)} / \|\tilde{\mathbf{w}}_i^{(n)}\|, l^{(n)} \right),$$

so $\mathbb{G}(\mathbf{x}^{(n+1)})\mathbb{U}_{\text{SD}}^{(n)}$ can be approximated by $\mathbb{Y}_{\text{SD}}^{(n)}$. With Rayleigh–Ritz methods, we calculate $\{\boldsymbol{\eta}_i^{(n)}\}_{i=1}^k$ as the eigenvectors corresponding to the smallest k eigenvalues of a $K \times K$ matrix $\mathbb{P}_{\text{SD}}^{(n)} = (\mathbb{U}_{\text{SD}}^{(n)})^\top \mathbb{Y}_{\text{SD}}^{(n)}$ and accordingly set $\{\mathbf{v}_i^{(n+1)}\}_{i=1}^k$ as the Ritz vectors $\{\mathbb{U}_{\text{SD}}^{(n)} \boldsymbol{\eta}_i^{(n)}\}_{i=1}^k$. Because of numerical errors, we use the eigenvectors of $(\mathbb{P}_{\text{SD}}^{(n)} + (\mathbb{P}_{\text{SD}}^{(n)})^\top)/2$ as the symmetric part of $\mathbb{P}_{\text{SD}}^{(n)}$ practically.

We can implement the LOBPCG method similarly. By defining

$$(40) \quad \mathbf{w}_i^{(n)} = \mathbf{v}_{i-k}^{(n-1)}, \quad i = k+1, \dots, 2k,$$

we calculate (36) and (39) for $i = 1, \dots, 2k$ and define two K -column matrices as ($k \leq K \leq 3k$)

$$(41) \quad \mathbb{U}_{\text{CG}}^{(n)} = \left[\mathbf{v}_1^{(n)}, \dots, \mathbf{v}_k^{(n)}, \tilde{\mathbf{w}}_i^{(n)} / \|\tilde{\mathbf{w}}_i^{(n)}\| : \|\tilde{\mathbf{w}}_i^{(n)}\| > \epsilon_w, i = 1, \dots, 2k \right],$$

$$(42) \quad \mathbb{Y}_{\text{CG}}^{(n)} = \left[\mathbf{u}_1^{(n)}, \dots, \mathbf{u}_k^{(n)}, \mathbf{y}_i^{(n)} : \|\tilde{\mathbf{w}}_i^{(n)}\| > \epsilon_w, i = 1, \dots, 2k \right]$$

instead of (37) and (38). Especially, if $\|\tilde{\mathbf{w}}_i^{(n)}\|$ is close to ϵ_w , a reorthogonalization step for $\mathbf{w}_i^{(n)}$ is required for numerical stability. By Rayleigh–Ritz methods, we get $\{\mathbf{v}_i^{(n+1)}\}_{i=1}^k$ accordingly.

Algorithm 2. HiOSD-LOBPCG (or LOBPSD) method for a k -saddle.

Input: $k \in \mathbb{N}$, $l^{(0)} > 0$, $\mathbf{x}^{(0)} \in \mathcal{H}$, $\{\mathbf{v}_i^{(0)}\}_{i=1}^k \subset \mathcal{H}$ satisfying $\langle \mathbf{v}_j^{(0)}, \mathbf{v}_i^{(0)} \rangle = \delta_{ij}$.

```

1: Set  $n = 0$ , compute  $\mathbf{f}^{(0)} = \mathbf{F}(\mathbf{x}^{(0)})$ ;
2: repeat
3:    $\mathbf{x}^{(n+1)} = \mathbf{x}^{(n)} + \beta^{(n)} \mathbf{g}^{(n)}$ ;
4:   Calculate  $\mathbb{U}^{(n)}$  as (41) (or (37));
5:   Calculate  $\mathbb{Y}^{(n)}$  as (42) (or (38));
6:    $\mathbb{P}^{(n)} = (\mathbb{U}^{(n)})^\top \mathbb{Y}^{(n)}$ ;
7:   Calculate  $\{\boldsymbol{\eta}_i^{(n)}\}_{i=1}^k$  as the  $k$  lowest eigenvectors of  $(\mathbb{P}^{(n)} + (\mathbb{P}^{(n)})^\top)/2$ ;
8:    $\mathbf{v}_i^{(n+1)} = \mathbb{U}^{(n)} \boldsymbol{\eta}_i^{(n)}$ ;
9:    $l^{(n+1)} = \max \{l^{(n)}/(1 + \beta^{(n)}), \varepsilon\}$ ;
10:   $\mathbf{f}^{(n+1)} = \mathbf{F}(\mathbf{x}^{(n+1)})$ ;
11:   $n := n + 1$ ;
12: until  $\|\mathbf{f}^{(n)}\| < \epsilon_F$ ;
Output:  $\mathbf{x}^{(n)}, \mathbf{v}_1^{(n)}, \dots, \mathbf{v}_k^{(n)}$ .

```

To summarize, we present Algorithm 2 as a discrete scheme, where $\mathbf{g}^{(n)}$ and $\mathbf{d}_i^{(n)}$ are defined the same as Algorithm 1, and $\tilde{\mathbf{w}}_i^{(n)}$ in step 4 and $\mathbf{y}_i^{(n)}$ in step 5 are defined as in (36) and (39).

In each iteration step, the LOBPCG (or LOBPSD) method needs another $4k$ (or $2k$) force evaluations. Supposing \mathbb{T} is the identity, the approximate eigenvectors determined in Algorithm 1 are also on the subspace $\mathcal{U}^{(n)}$ spanned by (41) (or (37)), so the LOBPCG method brings better approximate eigenvectors in each iteration step. Furthermore, the LOBPCG method ensures that $\mathbf{v}_i^{(n+1)}$ is a better approximation than $\mathbf{v}_i^{(n)}$ without a proper choice of step sizes $\gamma_i^{(n)}$. Consequently, the LOBPCG method improves accuracy and robustness at the sacrifice of computation costs.

2.4. Determination of the step size. We present several approaches to determine the step sizes in Algorithm 1 and Algorithm 2.

Explicit Euler scheme. The simplest way of choosing the step sizes is an explicit Euler scheme, $\beta^{(n)} = \gamma_i^{(n)} = \Delta t$, where Δt is a positive constant. We can accelerate the convergence by increasing Δt , but an overlarge Δt often leads to numerical instability.

Line search method. The line search method can be used to search the optimal step size in each iteration step. The step sizes of \mathbf{x} are calculated by minimizing a value function of $\|\mathbf{F}(\mathbf{x})\|^2$, which achieves its global minima at all critical points of $E(\mathbf{x})$. Both exact and inexact line search methods can be applied in practice. An upper bound τ should be set for $\beta^{(n)} \|\mathbf{g}^{(n)}\|$ to avoid overlarge jump between states, and a lower bound should be set for $\beta^{(n)}$ so that it is possible to escape from a neighborhood of other critical points, which is different from minimization problems.

BB gradient method. By following the same spirit of the OSD [43], we adopt the BB gradient method [3] to determine the step sizes. The BB method chooses $\mathbb{H}^{(n)} = \beta^{(n)} \mathbb{I}$ as an approximation of the inverse of Hessians and imposes some quasi-Newton properties.

Let $\mathbf{g}^{(n)}$ denote the gradient at $\mathbf{x}^{(n)}$ as Algorithm 1 and 2, and define $\Delta \mathbf{x}^{(n)} = \mathbf{x}^{(n)} - \mathbf{x}^{(n-1)}$ and $\Delta \mathbf{g}^{(n)} = \mathbf{g}^{(n)} - \mathbf{g}^{(n-1)}$. By solving optimization subproblems $\min_{\beta^{(n)}} \|\Delta \mathbf{x}^{(n)} - \beta^{(n)} \Delta \mathbf{g}^{(n)}\|$ and $\min_{\beta^{(n)}} \|\Delta \mathbf{x}^{(n)} / \beta^{(n)} - \Delta \mathbf{g}^{(n)}\|$, we get the BB1 step size

$$(43) \quad \beta_{\text{BB1}}^{(n)} = \frac{\langle \Delta \mathbf{x}^{(n)}, \Delta \mathbf{x}^{(n)} \rangle}{\langle \Delta \mathbf{x}^{(n)}, \Delta \mathbf{g}^{(n)} \rangle}$$

and the BB2 step size

$$(44) \quad \beta_{\text{BB2}}^{(n)} = \frac{\langle \Delta \mathbf{x}^{(n)}, \Delta \mathbf{g}^{(n)} \rangle}{\langle \Delta \mathbf{g}^{(n)}, \Delta \mathbf{g}^{(n)} \rangle}.$$

In practice, since $|\langle \Delta \mathbf{x}^{(n)}, \Delta \mathbf{g}^{(n)} \rangle|$ could be small, we apply the BB2 formula (44) to solve the optimization problems (1), and the step size of \mathbf{x} is determined as

$$(45) \quad \beta^{(n)} = \min \left\{ \frac{\tau}{\|\mathbf{g}^{(n)}\|}, \left| \frac{\langle \Delta \mathbf{x}^{(n)}, \Delta \mathbf{g}^{(n)} \rangle}{\langle \Delta \mathbf{g}^{(n)}, \Delta \mathbf{g}^{(n)} \rangle} \right| \right\}$$

with the same upper bound τ as the line search method, which is set as 0.5 in numerical experiments. The absolute values are taken in (45) as a precaution to avoid negative step sizes due to the negative eigenvalues of $\mathbb{G}(\mathbf{x})$. Similarly, the step sizes of \mathbf{v}_i in Algorithm 1 are determined as

$$(46) \quad \gamma_i^{(n)} = \left| \frac{\langle \Delta \mathbf{v}_i^{(n)}, \Delta \mathbf{d}_i^{(n)} \rangle}{\langle \Delta \mathbf{d}_i^{(n)}, \Delta \mathbf{d}_i^{(n)} \rangle} \right|,$$

where $\Delta \mathbf{v}_i^{(n)} = \mathbf{v}_i^{(n)} - \mathbf{v}_i^{(n-1)}$, $\Delta \mathbf{d}_i^{(n)} = \mathbf{d}_i^{(n)} - \mathbf{d}_i^{(n-1)}$, $i = 1, 2, \dots, k$.

3. Numerical examples. In this section, we show some numerical examples to illustrate the efficiency and reliability of the HiOSD method. The error is calculated as the gradient norm $\|\mathbf{F}(\mathbf{x})\|$ because the exact saddle points are often unavailable in practical cases. The initial directions are determined as the eigenvectors corresponding to the smallest k eigenvalues of the Hessian $\mathbb{G}(\mathbf{x}^{(0)})$ using the LOBPCG method and dimer approximations.

3.1. Modified Biggs EXP6 functions. We first test our algorithms on a toy model. Consider a six-dimensional Biggs EXP6 function [32],

$$(47) \quad B(\mathbf{x}) = \sum_{i=1}^6 \left(x_3 e^{-t_i x_1} - x_4 e^{-t_i x_2} + x_6 e^{-t_i x_5} - y_i \right)^2,$$

where $t_i = \frac{i}{10}$ and $y_i = e^{-t_i} - 5e^{-10t_i} + 3e^{-4t_i}$. Because of symmetry, $B(\mathbf{x})$ has many global minima, such as $\hat{\mathbf{x}} = (1, 10, 1, 5, 4, 3)^\top$. We modify (47) with bounded quadratic arctangent terms to construct $\hat{\mathbf{x}}$ as a k -saddle, obtaining modified Biggs EXP6 functions

$$(48) \quad B_k(\mathbf{x}) = B(\mathbf{x}) - \sum_{i=1}^k s_i \arctan^2(x_i - \hat{x}_i) + \sum_{i=k+1}^6 s_i \arctan^2(x_i - \hat{x}_i).$$

By choosing $\mathbf{s} = (4, 8, 16, 8, 4, 2)$, $\hat{\mathbf{x}}$ becomes a k -saddle of $B_k(\mathbf{x})$ for $k = 2, 3, 4, 5$.

For comparison, we implement the HiOSD method with the explicit Euler scheme, the HiOSD method with the BB gradient method, and the HiOSD-LOBPSD method with the BB gradient method. The HiOSD-LOBPCG method performs almost the

same as the HiOSD-LOBPSD method in this low-dimensional case because the subspaces (33) and (34) are both large enough to calculate the eigenvectors accurately. An initial point $\mathbf{x}^{(0)} = (0, 9, 1, 5, 4, 3)^\top$ achieves convergence to $\hat{\mathbf{x}}$ for all algorithms and $k = 2, 3, 4, 5$. However, $\nabla^2 B_k(\mathbf{x}^{(0)})$ has exactly $k - 2$ negative eigenvalues, so the local information of $\mathbf{x}^{(0)}$ differs widely from that of $\hat{\mathbf{x}}$. The error tolerance ϵ_F is 10^{-10} in this example.

Because a small Δt for the explicit Euler scheme certainly leads to a large number of iteration steps, we need to determine an “optimal” step size Δt , when the HiOSD method with the explicit Euler scheme converges in the fewest iteration steps. For cases of $k = 2, 3, 4, 5$, Δt is set as 0.06, 0.10, 0.12, and 0.14, respectively, and we apply these Δt to the explicit Euler scheme. For the BB gradient method, the first step size is simply set as 0.01. Results are shown in Table 1 and Figure 1. The explicit Euler scheme takes the most force evaluations even with the optimal step sizes, while the BB gradient method takes fewer iteration steps and less computation. As shown in Figure 1, the explicit Euler scheme has a linear convergent rate, while the BB gradient method has a faster linear convergent rate. With the BB gradient method, the HiOSD method costs a few more iteration steps than the HiOSD-LOBPSD method but fewer force evaluations and less CPU time. This suggests to us that though calculating unstable directions accurately can reduce the number of iterations, it may be unnecessary to calculate the unstable directions accurately when \mathbf{x} is far from $\hat{\mathbf{x}}$.

To test the sensitivity and robustness of the HiOSD method, we add some perturbations on the initial point $\mathbf{x}^{(0)}$. When using $\mathbf{x}^{(0)} = (0, 9, 1, 5, 4, 3)^\top \pm 0.2e_i$ (e_i is the i th unit coordinate vector), all HiOSD dynamics for $k = 2, 3, 4, 5$ can converge to the original saddle points. If we increase initial perturbations $\mathbf{x}^{(0)} = (0, 9, 1, 5, 4, 3)^\top \pm 0.5e_i$, most dynamics still converge, but three perturbations including $+e_5, -e_6$ and $+e_6$ ($k = 2$) fail to converge. If we use larger perturbations $\mathbf{x}^{(0)} = (0, 9, 1, 5, 4, 3)^\top \pm e_i$, more perturbations including $-e_1, -e_4, +e_5, -e_6$, and $+e_4$ ($k = 2$), $+e_6$ ($k = 2$) will fail to converge.

3.2. Degenerate cases. Although the HiOSD method is developed under a minimax framework, our algorithms can be applied to some degenerate cases as well. A well-known example of degenerate saddle points is the monkey saddle [30], which is

TABLE 1

Comparison of HiOSD with Euler, HiOSD with BB, and HiOSD-LOBPSD with BB for the modified Biggs EXP6 functions B_k of $k = 2, 3, 4, 5$. Euler: the explicit Euler scheme. BB: the BB gradient method. Iter: number of iteration steps till convergence. F eval: number of force evaluations till convergence.

k	Algorithm		$\ \mathbf{x} - \hat{\mathbf{x}}\ _2$	$\ \mathbf{F}(\mathbf{x})\ _2$	Iter	F eval	Time/s
2	HiOSD	Euler	3.2e-11	9.4e-11	119	596	0.030
	HiOSD	BB	3.6e-12	2.2e-11	38	191	0.009
	HiOSD-LOBPSD	BB	1.8e-11	8.1e-11	31	260	0.011
3	HiOSD	Euler	2.3e-11	7.3e-11	71	498	0.022
	HiOSD	BB	5.5e-13	3.2e-12	36	253	0.011
	HiOSD-LOBPSD	BB	6.2e-12	5.9e-11	32	373	0.014
4	HiOSD	Euler	2.8e-11	8.9e-11	57	514	0.021
	HiOSD	BB	8.4e-12	3.1e-11	34	307	0.013
	HiOSD-LOBPSD	BB	1.9e-11	6.8e-11	31	384	0.015
5	HiOSD	Euler	1.3e-11	9.7e-11	162	1783	0.096
	HiOSD	BB	3.5e-11	8.2e-11	44	485	0.020
	HiOSD-LOBPSD	BB	6.4e-12	1.6e-11	34	429	0.018

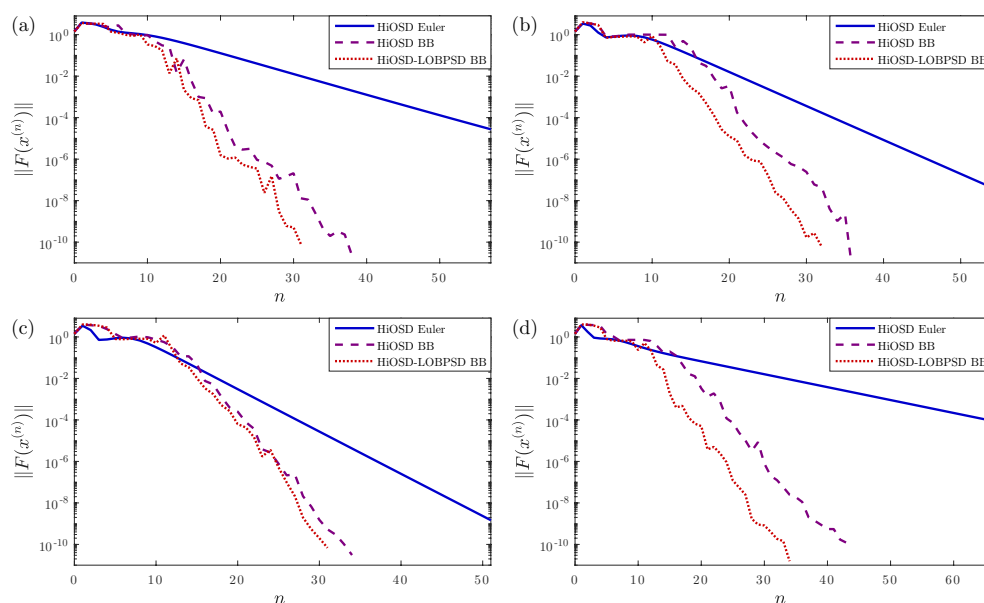


FIG. 1. Error reduction vs. iteration steps for HiOSD with Euler, HiOSD with BB, and HiOSD-LOBPSD with BB in the modified Biggs functions $B_k(\mathbf{x})$. (a) $k = 2$. (b) $k = 3$. (c) $k = 4$. (d) $k = 5$.

the origin of a two-dimensional function $E(x, y) = x^3 - 3xy^2$. It is obvious that this saddle point is not an extremum along any direction.

Now we consider computing some degenerate saddle points. The test function is a d -dimensional function ($d \geq 2$ is an even integer),

$$(49) \quad E_{d,p}(\mathbf{x}) = \frac{1}{p} \operatorname{Re} \sum_{j=1}^{d/2} \left(\sqrt[p]{j} \varphi(x_{2j-1}) + \frac{i}{\sqrt[p]{j}} \varphi(x_{2j}) \right)^p,$$

where $p \geq 3$ is an odd integer, i is the imaginary unit, and $\varphi(z) = \frac{z}{1+z^2}$. The origin is an isolated degenerate saddle point but is not a solution to any minimax problem (1). For instance, the landscape of function $E_{2,3}(\mathbf{x})$ is shown in Figure 2(a).

Nevertheless, we attempt to apply the HiOSD method with the explicit Euler scheme and the BB gradient method. We test our algorithms of $k = 3$ on $E_{6,3}$, and $k = 4$ on $E_{8,5}$, with an initial point $(0.4, \dots, 0.4)^\top$. Numerical results are shown in Figure 2(b, c) with the error tolerance $\epsilon_F = 10^{-12}$ and a step size $\Delta t = 1$. Until 10000 iteration steps, the HiOSD method with the explicit Euler scheme fails to achieve the error tolerance, while that with the BB gradient method succeeds within 40 iteration steps. It shows a sublinear convergent rate of the explicit Euler scheme and a linear convergent rate of the BB gradient method. In this degenerate case, the BB gradient method shows its advantage by accelerating the convergent rate from sublinear to linear.

3.3. Lane–Emden equation. Our HiOSD method can find multiple numerical solutions of some elliptic PDEs. Consider the Lane–Emden equation [27] with zero Dirichlet boundary conditions,

$$(50) \quad \begin{cases} -\Delta u(\mathbf{x}) = u^3(\mathbf{x}), & \mathbf{x} \in \Omega, \\ u(\mathbf{x}) = 0, & \mathbf{x} \in \partial\Omega, \end{cases}$$

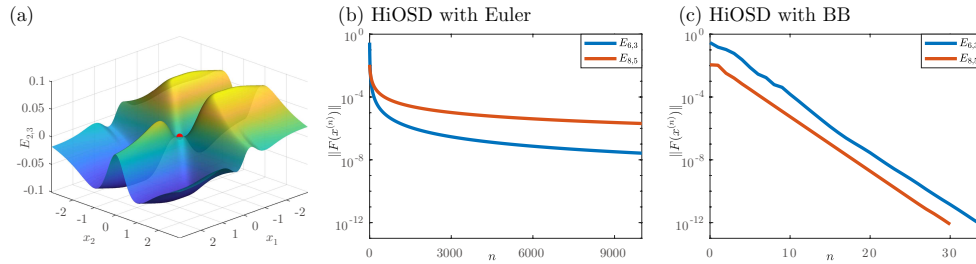


FIG. 2. (a) The landscape of $E_{2,3}$. The origin (red point) is a degenerate saddle point and is not a extremum along any direction. (b) Error reduction vs. iteration steps of the HiOSD method with the explicit Euler scheme for degenerate functions $E_{6,3}$ (blue line) and $E_{8,5}$ (orange line). (c) Error reduction vs. iteration steps of the HiOSD method with the BB gradient method for the degenerate functions.

where Ω is a bounded domain with Lipschitz boundaries. By the variational principle, solutions to the semilinear elliptic PDE (50) are critical points of the variational functional

$$(51) \quad E(u) = \int_{\Omega} \left(\frac{1}{2} |\nabla u(\mathbf{x})|^2 - \frac{1}{4} u^4(\mathbf{x}) \right) d\mathbf{x}, \quad u \in \mathcal{H}_0^1(\Omega).$$

The critical points that are not local extreme points can be identified as index- k saddle points.

We use an equivalent inner product of $\mathcal{H}_0^1(\Omega)$,

$$(52) \quad \langle u, v \rangle = \int_{\Omega} \nabla u(\mathbf{x}) \cdot \nabla v(\mathbf{x}) d\mathbf{x},$$

and the negative Fréchet derivative $F(u)$ with respect to the inner product (52) has the form $F(u) = -u + (-\Delta)^{-1}u^3$, where $(-\Delta)^{-1}u^3$ is the unique weak solution in $\mathcal{H}_0^1(\Omega)$ to a Poisson equation with respect to v ,

$$(53) \quad \begin{cases} -\Delta v(\mathbf{x}) = u^3(\mathbf{x}), & \mathbf{x} \in \Omega, \\ v(\mathbf{x}) = 0, & \mathbf{x} \in \partial\Omega, \end{cases}$$

The solutions to (50) represent density, so only positive solutions have physical significance. Nevertheless, we regard this as an example to test our algorithms, so the solutions with negative values are acceptable as well.

The homogeneous function $u_0(\mathbf{x}) \equiv 0$ is a local minimizer to (50). Alternatively, several “mound” functions are defined as

$$(54) \quad d_i(\mathbf{x}) = \begin{cases} \text{sign}(r_i) \left(1 + \cos \frac{\pi \|\mathbf{x} - \mathbf{x}_i\|_2}{\kappa r_i} \right), & \|\mathbf{x} - \mathbf{x}_i\|_2 \leq \kappa |r_i|, \\ 0 & \text{otherwise,} \end{cases}$$

where \mathbf{x}_i and r_i are determined according to the domain Ω . We require $\{\mathbf{x} : \|\mathbf{x} - \mathbf{x}_i\|_2 \leq |r_i|\} \subset \bar{\Omega}$, and $\kappa < 1$ is set as 0.9 to ensure that boundary conditions are satisfied. Then the initial guess of a k -saddle is constructed as $\sum_{i=1}^k d_i(\mathbf{x}) \|\nabla d_i\|_{L^2} / \|d_i^2\|_{L^2}$.

By using the finite element method, we approximate $\Omega \subset \mathbb{R}^2$ with a polygon and partition it into n_h triangles by MATLAB subroutine `initmesh` and `refinemesh`. We use the piecewise linear function space, and the Poisson equation (53) can be numerically solved by many efficient algorithms within reasonable costs, for example, the

MATLAB subroutine `assempte`. In the following cases, all step sizes are determined by the BB gradient method with the first step size $\Delta t = 0.02$, and the error tolerance ϵ_F is set as 10^{-5} .

Case 1. Ω is an asymmetric dumbbell-shaped domain consisting of a rectangle and two circles of radii 0.5 and 1, and a mesh of 5248 triangles is generated as shown in Figure 3(a). We attempt to find four index- k saddle points with conditions

- (A) $k = 2$, $\mathbf{x}_1 = (2, 0)$, $r_1 = 1$, $\mathbf{x}_2 = (-1, 0)$, $r_2 = 0.5$;
- (B) $k = 2$, $\mathbf{x}_1 = (2.5, 0)$, $r_1 = 0.5$, $\mathbf{x}_2 = (1.5, 0)$, $r_2 = -0.5$;
- (C) $k = 3$, $\mathbf{x}_1 = (2.5, 0)$, $r_1 = 0.5$, $\mathbf{x}_2 = (1.5, 0)$, $r_2 = -0.5$, $\mathbf{x}_3 = (-1, 0)$, $r_3 = 0.5$;
- (D) $k = 4$, $\mathbf{x}_1 = (2.5, 0)$, $r_1 = 0.5$, $\mathbf{x}_2 = (-1.25, 0)$, $r_2 = 0.25$, $\mathbf{x}_3 = (1.5, 0)$, $r_3 = -0.5$, $\mathbf{x}_4 = (-0.75, 0)$, $r_4 = -0.25$.

The numerical solutions of saddle points are shown in Figure 4, and the computation costs are shown in Table 2. The HiOSD method often takes more iteration steps but fewer force evaluations and less time because a Poisson equation is solved numerically to evaluating a force. Therefore, we apply the HiOSD method with the BB gradient method in the following cases. To test different initial guesses, we use $\sum_{i=1}^k 0.01d_i(\mathbf{x})$ with radii $r_i = \pm 0.1$, and the center of the mound \mathbf{x}_i is also perturbed for over

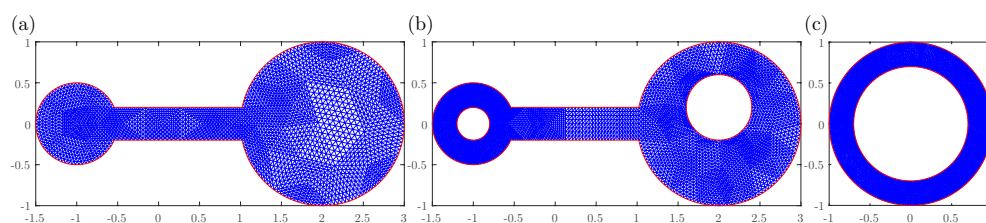


FIG. 3. The domains of the Lane–Emden equation and their triangular partitions. (a) A dumbbell-shaped domain (Case 1). (b) A dumbbell-shaped domain with two cavities inside (Case 2). (c) A concentric annular domain (Case 3).

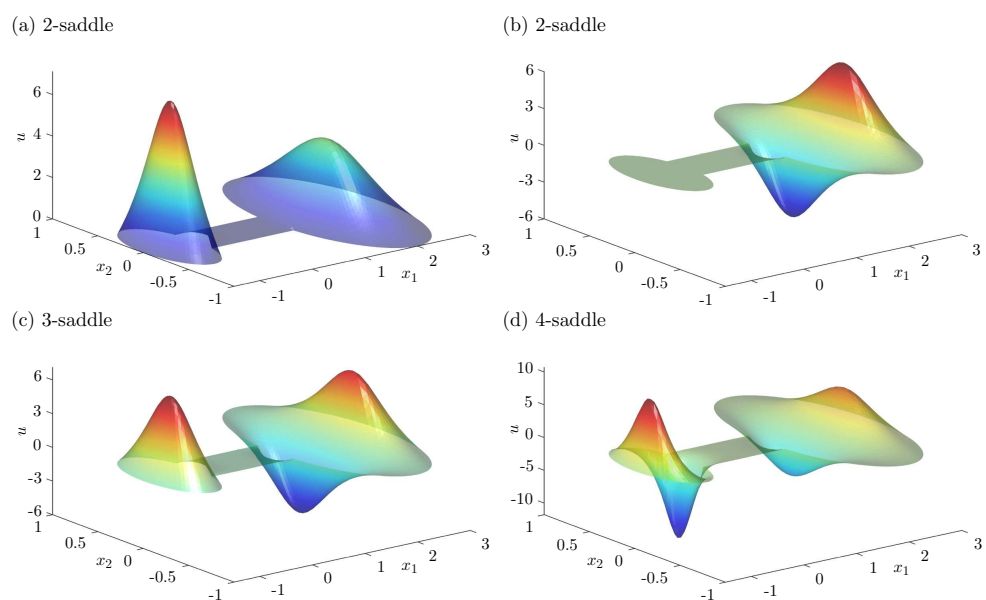


FIG. 4. Four saddle-point solutions to the Lane–Emden equation (Case 1).

TABLE 2
Comparison of the three algorithms for the Lane–Emden equation (Case 1).

	HiOSD			HiOSD-LOBPSD			HiOSD-LOBPCG		
	Iter	F_{eval}	Time/s	Iter	F_{eval}	Time/s	Iter	F_{eval}	Time/s
(A)	42	211	3.1	45	406	5.8	42	543	7.5
(B)	105	526	7.3	96	865	12.0	88	1139	15.7
(C)	102	715	10.1	86	1119	15.4	102	1933	26.7
(D)	169	1522	21.3	161	2738	37.9	157	3918	54.9

0.1 distance. The numerical tests show all perturbations can converge to the saddle points, indicating the robustness of the HiOSD method.

Case 2. Ω is the asymmetric dumbbell-shaped domain with two circular cavities of radii 0.4 and 0.2 inside, and a mesh of 9088 triangles is generated as shown in Figure 3(b). We attempt to find twelve index- k saddle points with conditions

- (A) $k = 2$, $\mathbf{x}_1 = (2, -0.6)$, $r_1 = 0.4$, $\mathbf{x}_2 = (-1.35, 0)$, $r_2 = 0.15$;
- (B) $k = 2$, $\mathbf{x}_1 = (2, -0.6)$, $r_1 = 0.4$, $\mathbf{x}_2 = (-0.5, 0)$, $r_2 = 0.2$;
- (C) $k = 2$, $\mathbf{x}_1 = (-0.5, 0)$, $r_1 = 0.2$, $\mathbf{x}_2 = (-1.35, 0)$, $r_2 = 0.15$;
- (D) $k = 2$, $\mathbf{x}_1 = (2.5, -0.4)$, $r_1 = 0.3$, $\mathbf{x}_2 = (1.5, -0.4)$, $r_2 = -0.3$;
- (E) $k = 2$, $\mathbf{x}_1 = (-1, 0.35)$, $r_1 = 0.15$, $\mathbf{x}_2 = (-1, -0.35)$, $r_2 = -0.15$;
- (F) $k = 3$, $\mathbf{x}_1 = (2, -0.6)$, $r_1 = 0.4$, $\mathbf{x}_2 = (-1.35, 0)$, $r_2 = 0.15$, $\mathbf{x}_3 = (-0.5, 0)$, $r_3 = 0.2$;
- (G) $k = 3$, $\mathbf{x}_1 = (2.5, -0.4)$, $r_1 = 0.3$, $\mathbf{x}_2 = (1.5, -0.4)$, $r_2 = -0.3$, $\mathbf{x}_3 = (-1.35, 0)$, $r_3 = 0.15$;
- (H) $k = 3$, $\mathbf{x}_1 = (2.5, -0.4)$, $r_1 = 0.3$, $\mathbf{x}_2 = (1.5, -0.4)$, $r_2 = -0.3$, $\mathbf{x}_3 = (-0.5, 0)$, $r_3 = 0.2$;
- (I) $k = 3$, $\mathbf{x}_1 = (-1, 0.35)$, $r_1 = 0.15$, $\mathbf{x}_2 = (-1, -0.35)$, $r_2 = -0.15$, $\mathbf{x}_3 = (2, -0.6)$, $r_3 = 0.4$;
- (J) $k = 4$, $\mathbf{x}_1 = (2.5, -0.4)$, $r_1 = 0.3$, $\mathbf{x}_2 = (1.5, -0.4)$, $r_2 = -0.3$, $\mathbf{x}_3 = (-1.35, 0)$, $r_3 = 0.15$, $\mathbf{x}_4 = (-0.5, 0)$, $r_4 = 0.2$;
- (K) $k = 4$, $\mathbf{x}_1 = (2.5, -0.4)$, $r_1 = 0.3$, $\mathbf{x}_2 = (-1, 0.35)$, $r_2 = 0.15$, $\mathbf{x}_3 = (1.5, -0.4)$, $r_3 = -0.3$, $\mathbf{x}_4 = (-1, -0.35)$, $r_4 = -0.15$;
- (L) $k = 4$, $\mathbf{x}_1 = (2.5, -0.4)$, $r_1 = -0.3$, $\mathbf{x}_2 = (1.5, -0.4)$, $r_2 = 0.3$, $\mathbf{x}_3 = (-1.35, 0)$, $r_3 = 0.15$, $\mathbf{x}_4 = (-0.5, 0)$, $r_4 = -0.2$.

Numerical solutions of various saddle points are shown in Figure 5.

Case 3. Ω is a concentric annular domain with an inner radius 0.7 and an outer radius 1. Since Ω is geometrically symmetric, any rotation around the origin of a solution is also a solution, so a nonradial solution is degenerated. By the finite element method, a mesh of 8192 triangles is generated as shown in Figure 3(c), where this symmetry is damaged. We attempt to find four index- k saddle points with conditions

- (A) $k = 2$, $\mathbf{x}_i = 0.85 \times (\cos \frac{2\pi i}{k}, \sin \frac{2\pi i}{k})$, $r_i = 0.15$;
- (B) $k = 3$, $\mathbf{x}_i = 0.85 \times (\cos \frac{2\pi i}{k}, \sin \frac{2\pi i}{k})$, $r_i = 0.15$;
- (C) $k = 4$, $\mathbf{x}_i = 0.85 \times (\cos \frac{2\pi i}{k}, \sin \frac{2\pi i}{k})$, $r_i = (-1)^i \times 0.15$;
- (D) $k = 6$, $\mathbf{x}_i = 0.85 \times (\cos \frac{2\pi i}{k}, \sin \frac{2\pi i}{k})$, $r_i = (-1)^i \times 0.15$.

Four saddle-point solutions are shown in Figure 6.

4. Discussions and conclusions. In this paper, we propose a HiOSD method for finding index- k saddle points without evaluations of Hessians. The HiOSD method can be regarded as a generalization of the OSD method for index-1 saddle points [43]. This saddle-searching problem is formulated by a minimax optimization problem, and a HiOSD dynamical system is presented to show its linearly stable steady state is

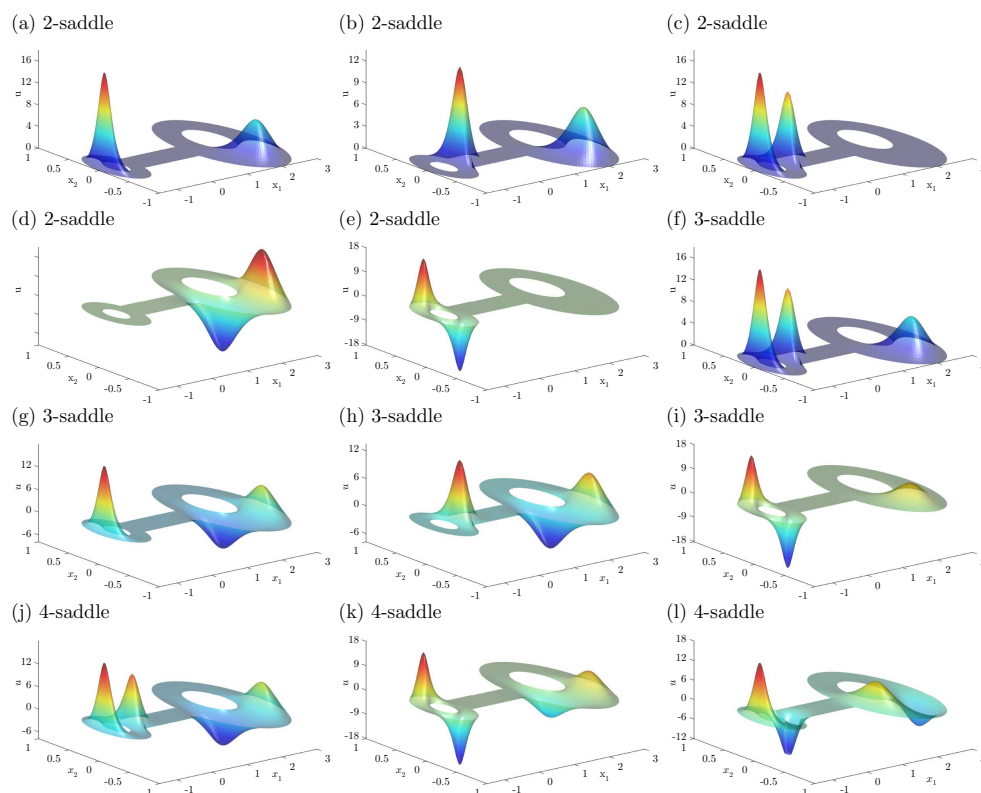


FIG. 5. Twelve saddle-point solutions to the Lane-Emden equation (Case 2).

exactly a k -saddle. We develop two algorithms as the implementation of HiOSD: one is to use simultaneous Rayleigh-quotient minimization technique, and the other is to apply the LOBPCG method in order to improve the accuracy and robustness. We also use various approaches to determine the step sizes, especially the BB gradient method. We test our algorithms by several examples to demonstrate the efficiency of the HiOSD method. By comparing the previous work by Quapp and Bofill in [33], we use dimer approximation to avoid the explicit calculation of Hessians, which is applicable to more practical problems. Under the optimization-based framework, we are able to apply efficient optimization methods to speed up the convergence and the linear stability theorem can be proved. The eigenvector calculation in the HiOSD method is a trade-off because we need to balance the computational cost and the accuracy of the eigenvector during the iterations. The adaptive strategy could be feasible to obtain the optimal convergence rate [16].

By using the HiOSD method, we provide a new approach to systematically explore the critical points of the energy functional, which can be labeled as index- k saddle points. Assuming that a k -saddle ($k \geq 1$) has been discovered, we can use the HiOSD algorithm to search the saddles with lower indices by following the unstable directions of this k -saddle, and repeat this searching process till reaching a minimizer. With this procedure, extensive critical points can be found systematically without initial guesses. We will pursue this idea of constructing the solution landscape that is a pathway map starting from a high-index saddle point and ending with minima elsewhere.

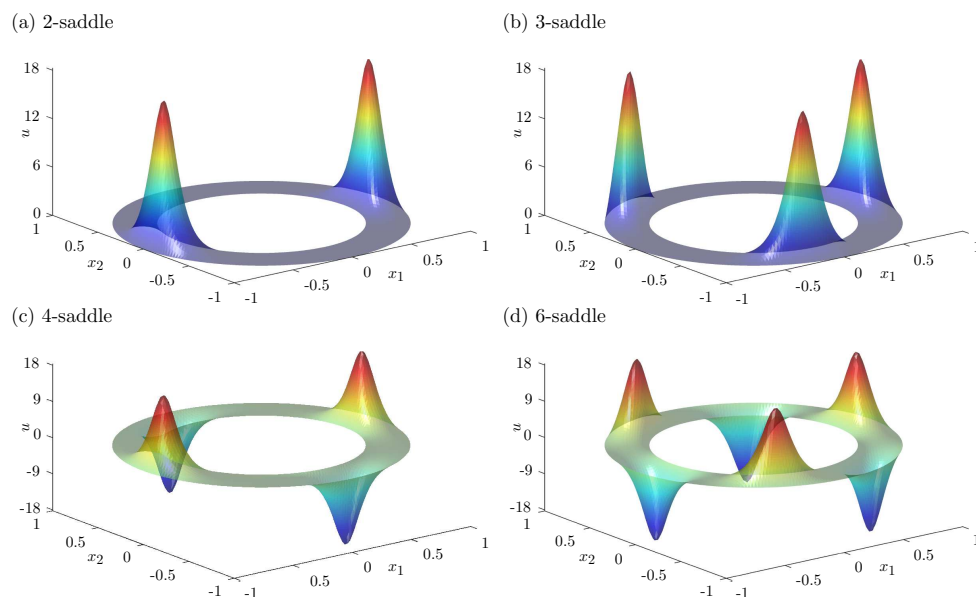


FIG. 6. Four saddle-point solutions to the Lane–Emden equation (Case 3).

The HiOSD method can be applied to solve many practical problems, for instance, finding the excited states for the Bose–Einstein condensation (BEC) [2]. A nonlinear Schrödinger equation, known as the Gross–Pitaevskii equation, can describe the dynamics of the BEC wave function. The stationary states of the BEC are the critical points of the Gross–Pitaevskii energy functional subject to a sphere constraint. Although the ground states of the BEC have been well studied [2], it is still not clear how to compute the excited states. We will develop the HiOSD method with a constraint to find the excited states of the BEC in future.

Acknowledgments. We would like to thank Professors Qiang Du, Zhijian Yang, and Weiqing Ren for fruitful discussions on the subject. We greatly appreciate Prof. Jianxin Zhou for providing their MATLAB codes of the Lane–Emden equation in [27]. We would also like to acknowledge the reviewers for their valuable comments.

REFERENCES

- [1] J. BAKER, *An algorithm for the location of transition states*, J. Comput. Chem., 7 (1986), pp. 385–395.
- [2] W. BAO AND Y. CAI, *Mathematical theory and numerical methods for Bose-Einstein condensation*, Kinet. Relat. Models, 6 (2013), pp. 1–135.
- [3] J. BARZILAI AND J. M. BORWEIN, *Two-point step size gradient methods*, IMA J. Numer. Anal., 8 (1988), pp. 141–148.
- [4] E. CANCÈS, F. LEGOLL, M.-C. MARINICA, K. MINOUKADEH, AND F. WILLAIME, *Some improvements of the activation-relaxation technique method for finding transition pathways on potential energy surfaces*, J. Chem. Phys., 130 (2009), 114711.
- [5] C. J. CERJAN AND W. H. MILLER, *On finding transition states*, J. Chem. Phys., 75 (1981), pp. 2800–2806.
- [6] X. CHENG, L. LIN, W. E, P. ZHANG, AND A.-C. SHI, *Nucleation of ordered phases in block copolymers*, Phys. Rev. Lett., 104 (2010), 148301.
- [7] G. M. CRIPPEN AND H. A. SCHERAGA, *Minimization of polypeptide energy: XI. The method of gentlest ascent*, Arch. Biochem. Biophys., 144 (1971), pp. 462–466.
- [8] Q. DU AND L. ZHANG, *A constrained string method and its numerical analysis*, Commun. Math. Sci., 7 (2009), pp. 1039–1051.

- [9] J. DUNCAN, Q. WU, K. PROMISLOW, AND G. HENKELMAN, *Biased gradient squared descent saddle point finding method*, J. Chem. Phys., 140 (2014), 194102.
- [10] W. E, W. REN, AND E. VANDEN-EIJNDEN, *String method for the study of rare events*, Phys. Rev. B, 66 (2002), 052301.
- [11] W. E, W. REN, AND E. VANDEN-EIJNDEN, *Simplified and improved string method for computing the minimum energy paths in barrier-crossing events*, J. Chem. Phys., 126 (2007), 164103.
- [12] W. E AND E. VANDEN-EIJNDEN, *Transition-path theory and path-finding algorithms for the study of rare events*, Annu. Rev. Phys. Chem., 61 (2010), pp. 391–420.
- [13] W. E AND X. ZHOU, *The gentlest ascent dynamics*, Nonlinearity, 24 (2011), pp. 1831–1842.
- [14] J. B. FORESMAN, M. HEAD-GORDON, J. A. POPL, AND M. J. FRISCH, *Toward a systematic molecular orbital theory for excited states*, J. Phys. Chem., 96 (1992), pp. 135–149.
- [15] W. GAO, J. LENG, AND X. ZHOU, *An iterative minimization formulation for saddle point search*, SIAM J. Numer. Anal., 53 (2015), pp. 1786–1805.
- [16] W. GAO, J. LENG, AND X. ZHOU, *Iterative minimization algorithm for efficient calculations of transition states*, J. Comput. Phys., 309 (2016), pp. 69–87.
- [17] N. GOULD, C. ORTNER, AND D. PACKWOOD, *A dimer-type saddle search algorithm with preconditioning and linesearch*, Math. Comp., 85 (2016), pp. 2939–2966.
- [18] W. J. GRANTHAM, *Gradient transformation trajectory following algorithms for determining stationary min-max saddle points*, in Advances in Dynamic Game Theory, Birkhäuser, Boston, MA, 2007, pp. 639–657.
- [19] S. GU AND X. ZHOU, *Convex splitting method for the calculation of transition states of energy functional*, J. Comput. Phys., 353 (2018), pp. 417–434.
- [20] D. HEIDRICH AND W. QUAPP, *Saddle points of index 2 on potential energy surfaces and their role in theoretical reactivity investigations*, Theor. Chim. Acta, 70 (1986), pp. 89–98.
- [21] G. HENKELMAN, G. JÓHANNESSON, AND H. JÓNSSON, *Methods for finding saddle points and minimum energy paths*, in Theoretical Methods in Condensed Phase Chemistry, Springer, Dordrecht, 2002, pp. 269–302.
- [22] G. HENKELMAN AND H. JÓNSSON, *A dimer method for finding saddle points on high dimensional potential surfaces using only first derivatives*, J. Chem. Phys., 111, pp. 7010–7022.
- [23] G. HENKELMAN AND H. JÓNSSON, *Improved tangent estimate in the nudged elastic band method for finding minimum energy paths and saddle points*, J. Chem. Phys., 113 (2000), pp. 9978–9985.
- [24] G. HENKELMAN, B. P. UBERUAGA, AND H. JÓNSSON, *A climbing image nudged elastic band method for finding saddle points and minimum energy paths*, J. Chem. Phys., 113 (2000), pp. 9901–9904.
- [25] A. HEYDEN, A. T. BELL, AND F. J. KEIL, *Efficient methods for finding transition states in chemical reactions: Comparison of improved dimer method and partitioned rational function optimization method*, J. Chem. Phys., 123 (2005), 224101.
- [26] A. V. KNYAZEV, *Toward the optimal preconditioned eigensolver: Locally optimal block preconditioned conjugate gradient method*, SIAM J. Sci. Comput., 23 (2001), pp. 517–541.
- [27] Y. LI AND J. ZHOU, *A minimax method for finding multiple critical points and its applications to semilinear PDEs*, SIAM J. Sci. Comput., 23 (2001), pp. 840–865.
- [28] D. E. LONGSINE AND S. F. MCCORMICK, *Simultaneous Rayleigh-quotient minimization methods for $Ax = \lambda Bx$* , Linear Algebra Appl., 34 (1980), pp. 195–234.
- [29] E. MACHADO-CHARRY, L. K. BÉLAND, D. CALISTE, L. GENOVESE, T. DEUTSCH, N. MOUSSEAU, AND P. POCHET, *Optimized energy landscape exploration using the ab initio based activation-relaxation technique*, J. Chem. Phys., 135 (2011), 034102.
- [30] J. MILNOR, *Morse Theory*, Princeton University Press, Princeton, NJ, 2016.
- [31] R. A. MIRON AND K. A. FICHTHORN, *The step and slide method for finding saddle points on multidimensional potential surfaces*, J. Chem. Phys., 115 (2001), pp. 8742–8747.
- [32] J. J. MORÉ, B. S. GARBOW, AND K. E. HILLSTROM, *Testing unconstrained optimization software*, ACM Trans. Math. Software, 7 (1981), pp. 17–41.
- [33] W. QUAPP AND J. M. BOFILL, *Locating saddle points of any index on potential energy surfaces by the generalized gentlest ascent dynamics*, Theor. Chem. Accounts, 133 (2014), 1510.
- [34] W. REN AND E. VANDEN-EIJNDEN, *A climbing string method for saddle point search*, J. Chem. Phys., 138 (2013), 134105.
- [35] A. SAMANTA, M. E. TUCKERMAN, T.-Q. YU, AND W. E, *Microscopic mechanisms of equilibrium melting of a solid*, Science, 346 (2014), pp. 729–732.
- [36] D. SHEPPARD, R. TERRELL, AND G. HENKELMAN, *Optimization methods for finding minimum energy paths*, J. Chem. Phys., 128 (2008), 134106.
- [37] D. WALES, *Energy Landscapes: Applications to Clusters, Biomolecules and Glasses*, Cambridge University Press, Cambridge, UK, 2003.

- [38] Y. WANG AND J. LI, *Phase field modeling of defects and deformation*, Acta Materialia, 58 (2010), pp. 1212–1235.
- [39] P. YU, Q. NIE, C. TANG, AND L. ZHANG, *Nanog induced intermediate state in regulating stem cell differentiation and reprogramming*, BMC Syst. Biol., 12 (2018), 22.
- [40] J. ZHANG AND Q. DU, *Shrinking dimer dynamics and its applications to saddle point search*, SIAM J. Numer. Anal., 50 (2012), pp. 1899–1921.
- [41] L. ZHANG, L.-Q. CHEN, AND Q. DU, *Morphology of critical nuclei in solid-state phase transformations*, Phys. Rev. Lett., 98 (2007), 265703.
- [42] L. ZHANG, L.-Q. CHEN, AND Q. DU, *Simultaneous prediction of morphologies of a critical nucleus and an equilibrium precipitate in solids*, Commun. Comput. Phys., 7 (2010), pp. 674–682.
- [43] L. ZHANG, Q. DU, AND Z. ZHENG, *Optimization-based shrinking dimer method for finding transition states*, SIAM J. Sci. Comput., 38 (2016), pp. A528–A544.
- [44] L. ZHANG, K. RADTKE, L. ZHENG, A. Q. CAI, T. F. SCHILLING, AND Q. NIE, *Noise drives sharpening of gene expression boundaries in the zebrafish hindbrain*, Mol. Syst. Biol., 8 (2012), 613.
- [45] L. ZHANG, W. REN, A. SAMANTA, AND Q. DU, *Recent developments in computational modelling of nucleation in phase transformations*, NPJ Comput. Materials, 2 (2016), 16003.